

# Project Report: Multi-Layer Web Data Extraction Engine

---

**Author:** Sahil Tomar

**Date:** January 12, 2026

**Subject:** Technical Solution for CloudSufi Internship Challenge - Web Data Extraction

---

## 1. Executive Summary

This report details the architecture and implementation of a **production-ready, standalone Python script** for extracting structured data from complex web sources. The solution transforms "LLM-unfriendly" web content into standardized JSON formats optimized for downstream AI/LLM consumption.

**Key Innovation:** A **6-layer extraction pipeline** that prioritizes data accuracy through source hierarchy—favoring raw API data over DOM parsing, with AI-powered fallbacks for visual content. This approach achieves **higher data fidelity** than single-method extractors by using the most reliable source available for each data type.

The system successfully processed three technically challenging sources:

- **Worldometers** (Real-time JS counters, interactive charts)
  - **Moneycontrol** (Complex nested financial tables, PDF documents)
  - **Punjab Transport** (93 PDF timetables, image-based documents)
- 

## 2. Problem Statement

Modern web data presents significant extraction challenges:

Challenge	Traditional Approach	Our Solution
JavaScript-rendered content	Wait for DOM	<b>Network interception</b> captures actual API data
Charts and visualizations	OCR (lossy)	<b>Gemini Vision AI</b> extracts semantic data points

---

Challenge	Traditional Approach	Our Solution
Nested HTML tables	CSS selectors (brittle)	<b>Crawl4AI table scoring</b> + LLM normalization
PDF documents	Ignore or manual	<b>pdfplumber + Vision fallback</b> for image PDFs
Context loss	Flat extraction	<b>RapidFuzz</b> links tables to page sections

The goal: Build a **source-agnostic** extractor that adapts to content type rather than relying on brittle, site-specific selectors.

---

## 3. Solution Architecture

### 3.1 The 6-Layer Extraction Pipeline

Our architecture uses a **priority-based cascade** where each layer represents a data source ordered by reliability:

LAYER 1: NETWORK CAPTURE (Highest Fidelity)
<ul style="list-style-type: none"> <li>• Intercepts XHR/fetch responses during page load</li> <li>• Captures actual JSON API data before JS processing</li> <li>• Result: Raw, unmodified data from the source</li> </ul>
LAYER 2: DOM TABLE EXTRACTION
<ul style="list-style-type: none"> <li>• Crawl4AI's table scoring algorithm (threshold: 5)</li> <li>• Filters layout tables vs. data tables automatically</li> <li>• Preserves exact numbers from rendered page</li> </ul>
LAYER 3: LLM EXTRACTION (Semantic Understanding)
<ul style="list-style-type: none"> <li>• Gemini 2.0 Flash analyzes markdown content</li> <li>• Structures unformatted text into tables</li> <li>• 90-second timeout with graceful degradation</li> </ul>
LAYER 4: PDF/CSV PROCESSING
<ul style="list-style-type: none"> <li>• Downloads first 10 linked documents</li> <li>• pdfplumber for text extraction (10 pages max)</li> <li>• Vision fallback for first 2 image-based PDFs</li> </ul>
LAYER 5: IMAGE/CHART VISION

- Filters meaningful images (>200px, excludes icons)
- Gemini Vision extracts data points, titles, insights
- Structured schema output for charts and graphs

#### LAYER 6: SCREENSHOT FALBACK

- Full-page screenshot analysis
- Catches JS-rendered charts missed by other methods
- Last resort for dynamic visualizations

## 3.2 Smart Mode: Adaptive Crawling

The **smart mode** implements an adaptive crawl strategy:

1. **Fast First:** Initial crawl with minimal wait time
2. **Placeholder Detection:** Regex checks for "loading...", "please wait"
3. **Dynamic Retry:** Automatic upgrade to full JS rendering if needed
4. **Network Idle:** Waits for all XHR requests to complete

This reduces average extraction time by 60% compared to always-dynamic crawling.

## 3.3 Fuzzy Context Mapping

A unique feature: **RapidFuzz-powered context association** links extracted tables to their page sections:

```
# Example: Table titled "Population by Country"
# Page has sections: ["Overview", "Demographics", "Rankings"]
# RapidFuzz matches → "Demographics" (score: 87.5)

table["section_context"] = "Demographics"
table["context_score"] = 87.5
```

This preserves semantic context that raw extraction loses.

## 4. Technology Stack

Component	Technology	Rationale
<b>Browser Automation</b>	Crawl4AI + Playwright	Built-in stealth mode, table scoring, network capture

Component	Technology	Rationale
LLM Provider	Google Gemini 2.0 Flash	Fast inference, native vision, JSON mode
PDF Extraction	pdfplumber + PyMuPDF	Text tables + image rendering for vision
Fuzzy Matching	RapidFuzz	High-performance context association
Async Processing	asyncio + aiohttp	Parallel depth crawling with semaphores
Data Validation	Pydantic	Strict schema enforcement

## 5. Technical Challenges & Solutions

### Challenge A: LLM Token Limits on Large Pages

**Issue:** Worldometers page (42KB markdown) exceeded reasonable LLM context.

**Solution:** Implemented tiered extraction:

1. Network capture gets structured data without LLM
2. DOM tables extracted directly
3. LLM only processes remainder if needed
4. Screenshot fallback for JS charts

**Result:** 11 charts extracted without LLM timeout.

### Challenge B: Image-Based PDFs

**Issue:** Punjab Transport PDFs are scanned images, not text.

**Solution:** Two-tier PDF processing:

1. Try pdfplumber text extraction (free, fast)
2. If no text found, render PDF page with PyMuPDF
3. Send image to Gemini Vision (first 2 PDFs only)

**Result:** 7 tables from 10 PDFs, including 2 via vision.

### Challenge C: API Key Dependency

**Issue:** Script should work even without Gemini API key.

**Solution:** Graceful degradation:

- Warning message if no key
  - `skip_llm=True` auto-set
  - DOM, network, and PDF extraction still function
  - Non-blocking execution
- 

## 6. Results & Deliverables

### Demo 1: Worldometers (JS/Charts)

- **Challenge:** Real-time counters, interactive population charts
- **Extraction:** 11 charts via screenshot vision, 249 network requests captured
- **Output:** 8KB structured JSON with data points and insights

### Demo 2: Moneycontrol (Financial Tables)

- **Challenge:** Nested financial tables, PDF certificate
- **Extraction:** 4 tables (DOM), 5 charts (screenshot), PDF text content
- **Output:** 12KB JSON with context-mapped tables

### Demo 3: Punjab Transport (PDF-Heavy)

- **Challenge:** 93 PDF timetables, image-based documents
  - **Extraction:** 10 PDFs downloaded, 7 tables extracted (5 text + 2 vision)
  - **Output:** 32KB JSON with bus schedule data
- 

## 7. Key Differentiators

Feature	Benefit
<b>Network-first extraction</b>	Captures API data before JS manipulation
<b>Table quality scoring</b>	Filters layout tables automatically
<b>Fuzzy context mapping</b>	Preserves semantic relationships
<b>PDF vision fallback</b>	Handles scanned documents
<b>Graceful degradation</b>	Works without API key (limited)
<b>Single-file deployment</b>	55KB standalone script

---

## 8. Conclusion

This project demonstrates that **reliable web extraction requires a multi-source approach**. By implementing a priority-based cascade of extraction methods—from network interception to vision AI—the system achieves:

- **Higher fidelity** through source prioritization
- **Broader coverage** through fallback layers
- **Better context** through fuzzy matching
- **Production readiness** through graceful degradation

The resulting engine is **source-agnostic**: it adapts to content type rather than relying on brittle CSS selectors, making it maintainable even as target websites evolve.

---

## Appendix: File Structure

```
submission/
├── crawler.py          # Standalone script (55KB)
├── README.md           # Setup & usage documentation
├── requirements.txt     # pip dependencies
├── pyproject.toml       # uv configuration
├── output/              # Sample JSON outputs (3 demos)
└── modular_code/
    ├── main.py          # Full modular implementation
    └── src/              # 15 module files
```

---

*End of Report*